



PATENT ABSTRACTS OF JAPAN

(11) Publication number: **06223121 A**

(43) Date of publication of application: **12.08.94**

(51) Int. Cl.

G06F 15/40
G06F 15/40
G06K 9/72

(21) Application number: **05008734**

(71) Applicant: **NEC CORP**

(22) Date of filing: **22.01.93**

(72) Inventor: **KANEDA SATORU**

(54) **INFORMATION RETRIEVING DEVICE**

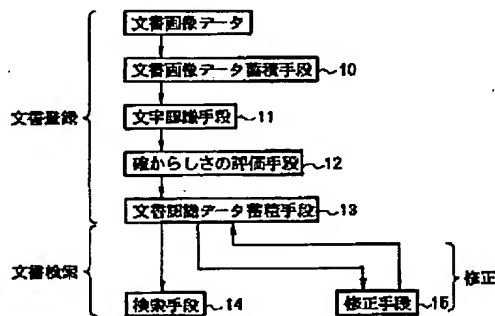
retrieval omission can be prevented even without performing correction.

(57) Abstract:

COPYRIGHT: (C)1994,JPO&Japio

PURPOSE: To prevent retrieval omission caused by uncertain recognition and to perform correction later in the information retrieving device and an information correcting device for performing retrieval and correction with the character recognized result as an object.

CONSTITUTION: A document image is read from a document image data storage means 10, and a character recognizing means 11 divides a character image at every character, performs character recognition and outputs candidate characters and the certainly A evaluating means 12 for certainty selects any candidate so as to prevent retrieval omission from the candidate characters and the certainty, enumerates the candidate and preserves it in list in a document recognition data storage means 13. A retrieving means 14 performs collation for the unit of character on the condition that respective characters are contained in the candidate when reading document recognition data and performing character string retrieval. Therefore, when proper characters are contained in the candidate, the



(19) 日本国特許庁 (J P)

(12) 特 許 公 報 (B 2)

(11) 特許番号

第2586372号

(45) 発行日 平成 9 年 (1997) 2 月 26 日

(24) 登録日 平成 8 年 (1996) 12 月 5 日

(51) Int.Cl. ⁸	識別記号	庁内整理番号	F I	技術表示箇所
G 0 6 F 17/30			G 0 6 F 15/401	3 1 0 A
G 0 6 K 9/00		9061-5H	G 0 6 K 9/00	S
			G 0 6 F 15/40	3 7 0 A
				3 7 0 B
			15/403	3 1 0 C
請求項の数 3 (全 8 頁) 最終頁に続く				

(21) 出願番号 特願平5-8734

(22) 出願日 平成 5 年 (1993) 1 月 22 日

(65) 公開番号 特開平6-223121

(43) 公開日 平成 6 年 (1994) 8 月 12 日

(73) 特許権者 000004237

日本電気株式会社

東京都港区芝五丁目 7 番 1 号

(72) 発明者 金田 悟

東京都港区芝五丁目 7 番 1 号 日本電気株式会社内

(74) 代理人 弁理士 京本 直樹 (外 2 名)

審査官 ▲吉▼田 耕一

(56) 参考文献 特開 平 3 - 50692 (J P, A)

特開 平 2 - 82380 (J P, A)

特開 平 3 - 160580 (J P, A)

特開 昭 62 - 44878 (J P, A)

特開 平 4 - 37971 (J P, A)

(54) 【発明の名称】 情報検索装置及び情報検索方法

1

(57) 【特許請求の範囲】

【請求項 1】 文書画像データを文字認識して得られた文書認識データ群から、任意の検索キーワードを含む文書認識データを検索する情報検索装置において、
入力された文書画像データを蓄積する文書画像データ蓄積手段と、
文書画像データに含まれる文字部分の文字パターンを認識し、候補となる文字コードを選択し、文字コードの確からしさの推定値を求める文字認識手段と、
文字コードの確からしさの推定値の和が、複数の文字コードを格納する文字コードリスト中に正確な文字が含まれるように定めた所定の確率を越えるまで、前記推定値の高い順に候補となる文字コードを前記文字コードリストに追加し、前記文字コードリストの中の候補となる文字コードの数が 1 つの場合は、先頭の候補となる文字コ

2

ードを選択し、前記リスト中の文字コードの数が定められたしきい値以内の場合は、候補となる文字コードが複数あることを示す認識コードをとまなう先頭の候補となる文字コードを含む複数の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値を越えた場合は、候補となる文字コードが多数あることを示す識別コードを選択する確からしさの評価手段と、
これら選択されたコードを文書認識データとして蓄積する文書認識データ蓄積手段と、
入力手段から入力した任意のキーワードを前記文書認識データ蓄積手段から検索する検索手段を備えることを特徴とする情報検索装置。

【請求項 2】 前記文書認識データ蓄積手段の文書認識データの中の複数の候補となる文字コードから正しい文字

コードの選択や、前記文字コードが多数あることを示す識別コードに正しい文字コードを入力する修正手段を備えることを特徴とする請求項1記載の情報検索装置。

【請求項3】文書画像データを文字認識して得られた文書認識データ群から、任意の検索キーワードを含む文書認識データを検索する情報検索方法において、

入力された文書画像データを蓄積する文書画像データ蓄積ステップと、

文書画像データに含まれる文字部分の文字パターンを認識し、候補となる文字コードを選択し、文字コードの確からしさの推定値を求める文字認識ステップと、

文字コードの確からしさの推定値の和が、複数の文字コードを格納する文字コードリスト中に正確な文字が含まれるように定めた所定の確率を越えるまで、推定値の高い順に候補となる文字コードを前記文字コードリストに追加し、前記文字コードリストの中の候補となる文字コードの数が1つの場合は、先頭の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値以内の場合は、候補となる文字コードが複数あることを示す認識コードをとともう先頭の候補となる文字コードを含む複数の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値を越えた場合は、候補となる文字コードが多数あることを示す識別コードを選択する確からしさの評価ステップと、

これら選択されたコードを文書認識データとして蓄積する文章認識データ蓄積ステップと、

入力した任意のキーワードを、蓄積された文書認識データから検索する検索ステップを備えることを特徴とする情報検索方法。

【発明の詳細な説明】

【0001】

【産業上の利用分野】本発明は文字データの検索を行う情報検索装置および情報検索方法に関し、特に文書画像から得られた文字データの情報検索装置及び情報検索方法に関する。

【0002】

【従来の技術】近年、文書情報の増大にともない、文書画像を電子化して蓄積（ファイル）しておき、検索して取り出す装置の需要が高まっている。このような電子ファイル装置において、文字認識を利用して文書画像の文書の内容に基づく検索を可能にする情報検索装置が望まれている。

【0003】

以下、上述した従来の情報検索装置について図を用いて説明する。図5は従来の文字認識を利用した情報検索装置のフローチャートである。

【0004】文書を登録する際は、文書をスキャナなどの画像入力手段1を利用して文書画像データとして入力し、文書画像蓄積手段2に蓄積する。さらに、文字認識手段3を利用して文書画像データに含まれている文字パ

タン部分を文字認識する。文字認識で認識が不確かな箇所が候補の文字が複数ある箇所については、キーボードなどで登録者が修正作業を行う。この認識結果データを文書画像データと対応づけて認識結果データ蓄積手段6に蓄積する。

【0005】前記のようにして登録した文書から目的の文書を検索する際は、キーボード7などからキーワードなどの検索条件を入力し、条件を満たす認識結果データを文字検索手段8などにより検索し、認識結果データに対応する文書画像データを出力する。

【0006】しかしながら、上記のような方式では、文書登録時に上記文字認識後に上記修正作業が必要になり、手間がかかる。

【0007】また、上記修正作業を行わないと、認識結果データに誤りが含まれる可能性があり、検索対象キーワードを文字列検索するとき、前記キーワードを含む認識結果データに誤りがあると、一致しないため検索漏れが生じる。検索漏れを防ぐために、キーワードの検索時に数文字までの不一致を許す照合手法や、検索文字列と類似したパタン形状の文字列を検索文字列と一緒に検索する手法が考えられてきた。こうした従来例として、信学技報CA87-25（1987年5月29日）、特開平4-158478号公報等に記載がある。

【0008】

【発明が解決しようとする課題】しかし、従来技術で行われていた、検索漏れを防ぐために数文字までの不一致を許す照合手法を使うと、不適当な検索結果が生じる（過検索）という問題がある。例えば、1文字までの誤りを許す検索手法で、検索文字列（“自由”など）を検索すると、まったく別の文字列（“理由”や“自然”など）とも一致していると判断してしまう。

【0009】また、検索文字列と類似形状の文字列と一緒に検索する手法では、文書画像に書体が異なる文字が含まれている場合など、誤認識の傾向があらかじめ予想されたものと異なる文字が含まれている場合に、検索漏れが起こる。

【0010】また、従来技術で、検索対象とする文書認識データに認識が不確かであった箇所や認識時の候補が何であったかといった情報が含まれていない場合、このデータだけを用いて修正作業するのは困難である。上記の文書認識データに含まれない情報は、別に保存しておく必要があり、管理が困難である。

【0011】

【課題を解決するための手段】図1は本発明の構成を示すブロック図である。図1に示すように、上記の課題を解決する第1の情報検索装置は、文書画像データを文字認識して得られた文書認識データ群から、任意の検索キーワードを含む文書認識データを検索する情報検索装置において、入力された文書画像データを蓄積する文書画像データ蓄積手段10と、文書画像データに含まれる文

字部分の文字パターンを認識し、候補となる文字コードを選択し、文字コードの確からしさの推定値を求める文字認識手段11と、文字コードの確からしさの推定値の和が、複数の文字コードを格納する文字コードリスト中に正確な文字が含まれるように定めた所定の確率を越えるまで、前記推定値の高い順に候補となる文字コードを前記文字コードリストに追加し、前記文字コードリストの中の候補となる文字コードの数が1つの場合は、先頭の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値以内の場合は、候補となる文字コードが複数あることを示す認識コードをとともなう先頭の候補となる文字コードを含む複数の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値を越えた場合は、候補となる文字コードが多数あることを示す識別コードを選択する確からしさの評価手段12と、これら選択されたコードを文書認識データとして蓄積する文書認識データ蓄積手段13と、入力手段から入力した任意のキーワードを前記文書認識データ蓄積手段から検索する検索手段14を備えることを特徴とする。第2の情報検索装置は第1の情報検索装置に加え、前記文書認識データ蓄積手段の文書認識データの中の複数の候補となる文字コードから正しい文字コードの選択や、前記文字コードが多数あることを示す識別コードに正しい文字コードを入力する修正手段15を備えることを特徴とする。

【0012】ここで、確からしさの評価手段12は、認識が不確かな文字については、候補となる複数の文字コードを列挙して出力する。ただし、候補の数が多い場合は、候補を列挙すると文書認識データが大きくなってしまったため、代わりに、全ての文字が候補であること（候補多数）を示す識別コードだけを出力する。

【0013】本発明の情報検索方法は、文書画像データを文字認識して得られた文書認識データ群から、任意の検索キーワードを含む文書認識データを検索する情報検索方法において、入力された文書画像データを蓄積する文書画像データ蓄積ステップと、文書画像データに含まれる文字部分の文字パターンを認識し、候補となる文字コードを選択し、文字コードの確からしさの推定値を求める文字認識ステップと、文字コードの確からしさの推定値の和が、複数の文字コードを格納する文字コードリスト中に正確な文字が含まれるように定めた所定の確率を越えるまで、推定値の高い順に候補となる文字コードを前記文字コードリストに追加し、前記文字コードリストの中の候補となる文字コードの数が1つの場合は、先頭の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値以内の場合は、候補となる文字コードが複数あることを示す認識コードをとともなう先頭の候補となる文字コードを含む複数の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値を越えた場合は、候補となる文

字コードが多数あることを示す識別コードを選択する確からしさの評価ステップと、これら選択されたコードを文書認識データとして蓄積する文章認識データ蓄積ステップと、入力した任意のキーワードを、蓄積された文書認識データから検索する検索ステップを備えることを特徴とする。

【0014】

【作用】確からしさの評価手段12は、文字認識手段11により文字認識時に得られた文字コードの確からしさの推定値の和が、複数の文字コードを格納する文字コードリスト中に正確な文字が含まれるように定めた所定の確率を越えるまで、推定値の高い順に候補となる文字コードを前記文字コードリストに追加し、前記文字コードリストの中の候補となる文字コードの数が1つの場合は、先頭の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値以内の場合は、候補となる文字コードが複数あることを示す認識コードをとともなう先頭の候補となる文字コードを含む複数の候補となる文字コードを選択し、前記リスト中の文字コードの数が定められたしきい値を越えた場合は候補となる文字コードが多数あることを示す識別コードを選択する。

【0015】このため、上記のようにして得られた文書認識データと検索対象キーワードとの文字ごとの照合をする際に、一致条件をキーワードの各文字が文書認識データの候補の中のどれかと一致すればよいとしておくことで、検索漏れを低減することができる。また、画像の条件が良い場合など高い確からしさで認識できる場合は候補を列挙しないので、余分な一致の発生が低減される。

【0016】また、文字認識の不確かさを評価して、誤りのありそうな箇所とその候補を選んで選択しているため、余分な候補が除去され、文書認識のデータ量が過大になるのが防がれる。

【0017】また、文書認識データに、認識が不確かな箇所と候補の情報が含まれるため、これを使って後から修正作業を行うことも可能である。

【0018】

【実施例】以下この発明の実施例について図面を参照しながら説明するが、この発明は以下の実施例に限定されるものではない。

【0019】図2は本発明の実施例の構成のブロック図を示すものである。図2において文書をスキャナ16などの画像読み取り手段で、画像データとして読み込み、文書画像データ蓄積手段17に蓄積する。文字認識手段18は、文書画像データ蓄積手段17から画像データを読み出して、文字が書かれた領域を識別して、文字パターンを認識し、候補となる文字コードとその確からしさを推定して確率値で出力する。

【0020】確からしさを推定するためには、パターンを

処理して得られるいくつかの統計量を軸とするベクトル空間を仮定して、あらかじめいろいろな文字についてその空間での座標を求めておく。そして、文字認識時に認識対象とする文字パターンを、同様に処理して上記ベクトル空間での座標を求め、学習時に近傍にあった座標を見つけ、それらとの距離から推定する。

【0021】確からしさの評価手段19は、図3に示すようなアルゴリズムに従う。各文字ごとに、文字認識手段が出力した候補のうち最も有力な候補をリストの先頭に入れる。そして、認識の確からしさをもとに、正しい文字がリストに含まれる確率Pを求める。確率Pが、あるしきい値 P_{th} より小さいならば、候補から次に有力な候補をリストに追加して、確率Pを求めなおして繰り返す。選出されたリストの中の候補の数（リストの長さ）が1個ならば、この文字列コードだけを出力する。リストの中の候補の数があるしきい値の個数 N_{th} より少ないなら、候補が複数個あることを示す識別コードと共に出力候補の文字コードを出力する。出力候補の数が N_{th} を越えていたならば、出力候補の文字コードの代わりに、全ての文字が候補であることを示す識別コード、すなわち、候補多数を示す識別コードを出力する。

【0022】識別コードの具体的な例としては、正規表現に準拠した記述が考えられる。例えば、候補として{"B"、"E"、"3"}を列挙する表現は、"

[BE3]"となり、認識不可の表現は、"."となる。NECというパターンを認識した結果、Nをはっきり認識し、Eの候補が{"B"、"E"、"3"}で、Cが候補多数であった場合、正規表現で記述すると、"N[BE3]*"となる。

【0023】文書認識データ蓄積手段20は、評価手段19が出力する文字コードと識別コードからなるコード列を文字認識データとして蓄積する。

【0024】文字列検索手段21は、キーボード24などの入力手段から入力した検索対象キーワードと文書認識データ蓄積手段20から読み出した文書認識データとを比較照合し、検索対象キーワードを含む文書認識データを検索する。

【0025】この文字列照合の例として図4に示したものは、文書認識データのテキストとキーワードを1文字1文字比較して、キーワードの文字列とテキストが全ての文字で一致している場合、キーワードの一致が成立したと判断するものである。ただし、候補複数を示す識別コードがあった場合は、該当するキーワードの文字がここに列挙された候補に含まれてたならば、この文字は一致しているとみなす。また、候補多数を示す識別コードがあった場合は、該当するキーワードの文字が何であっても、この文字は一致しているとみなす。

【0026】検索結果を知らせるために、上記の検索文字列が含まれる文書認識データ、あるいは、この文書認識データに対応する文書画像データをディスプレイ22

から表示する。

【0027】文書認識データ修正手段23は、文書認識データをディスプレイ22に表示し、複数の候補が列挙されている箇所については、ユーザーにキーボード24から正しいものを選択される。また、候補が特定されなかった箇所については、正しい文字コードをキーボード24から入力させる。

【0028】

【発明の効果】以上の実施例によれば、第1に、文字パターンの認識時に1つの候補だけでは確からしさを保証できない箇所は、確からしさの推定値が十分になるように複数の文字を候補にしたり、全ての文字を候補として指定するため、認識結果のテキストデータの大きさを過大にすることを防ぐことができる。

【0029】第2に、認識文字毎の認識の確からしさに応じて候補を列挙するため、検索漏れと過検索の低減を両立することができる。

【0030】第3に、文書画像に部分的なノイズや異フォントが含まれる場合など、誤認識の傾向が異なるときでも、個別に誤認識の傾向を示すデータ等を用意しなくても検索できる。

【0031】第4に、修正作業を後から行うことができるため、検索前の修正作業の手間と時間を省くことができる。

【0032】第5に、これらの修正作業を形態素解析などを行って、候補となる文字を選択した場合、この単語が辞書に含まれるかどうかによって候補を選択するような処理を行う場合も、これらの処理は処理装置に余裕ができるまで後回しできる。

【0033】第6に、認識結果内に複数の候補を記述する際の仕様が統一されてさえいれば、文字認識手段は異なってもかまわないため、英文や手書き文など文書ごとの特性に応じた異なる文字認識手段が作成した認識結果も同一の検索手段で検索することができる。

【0034】第7に、将来、より高性能な文書認識手段に切り替えた場合でも、それまでに蓄積した文書認識データを継続して利用することができる。

【図面の簡単な説明】

【図1】この発明の構成を示すブロック図である。

【図2】この発明の実施例を示すブロック図である。

【図3】確からしさの評価手段の実施例の動作を示すフローチャートである。

【図4】文字列検索の動作例を示す図である。

【図5】従来技術を示すブロック図である。

【符号の説明】

16 スキャナ

17 文書画像データ蓄積手段

18 文字認識手段

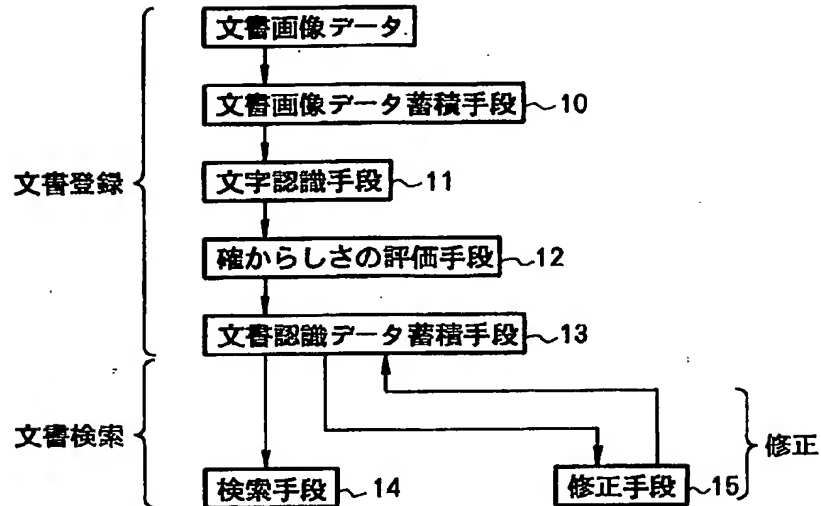
19 確からしさの評価手段

20 文書認識データ蓄積手段

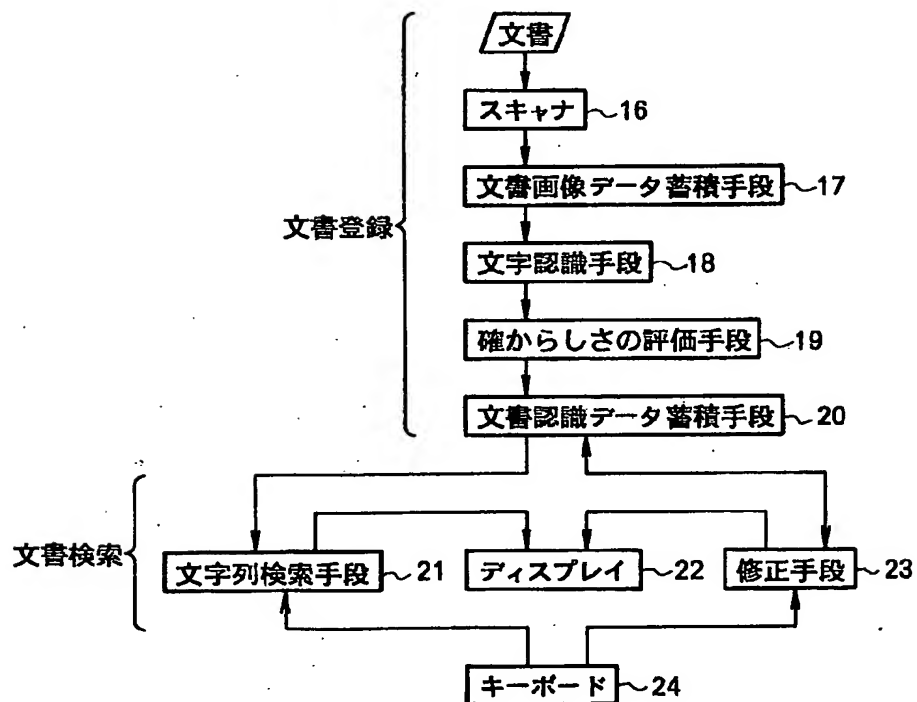
2.1 文字列検索手段

2.3 修正手段

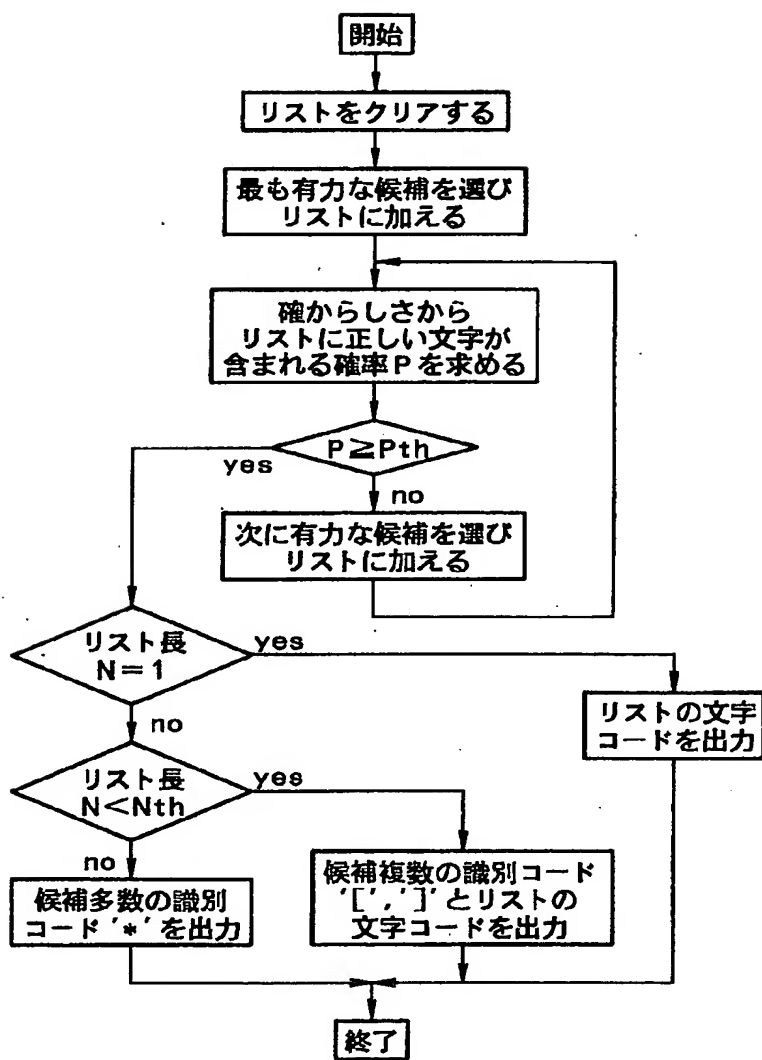
【図1】



【図2】



【図3】



【図4】

検索文字列: "NEC"

(1) 一致しない場合

文書認識データ: a b c d e f g N E T o p q r s t u
 文字の一致 : ↑↑↑
 O O x ... 一致不成立
 (O...一致、x...不一致)

(2) 一致する場合

文書認識データ: a b c d e f g N E C o p q r s t u
 文字の一致 : ↑↑↑
 O O O ... 一致成立

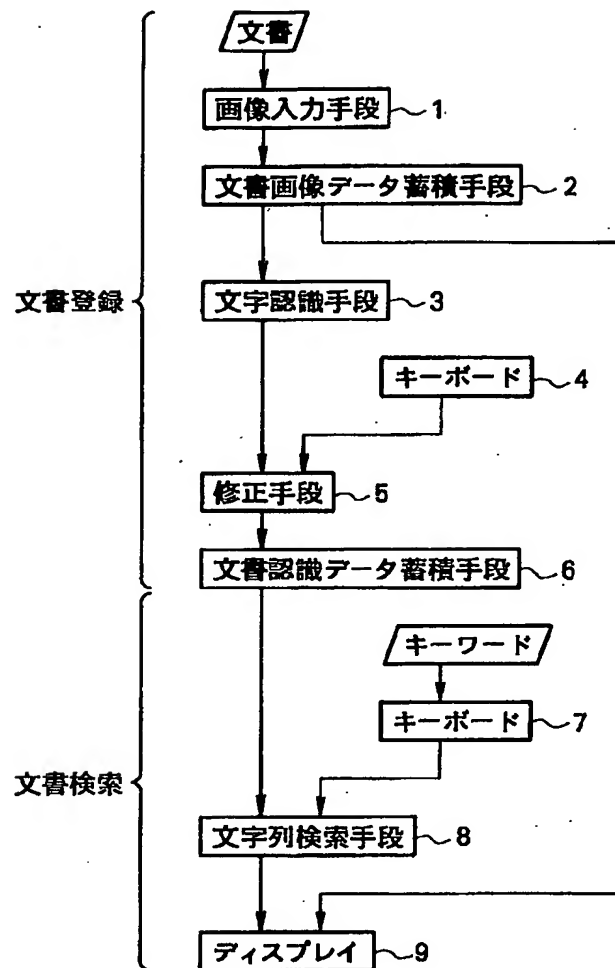
(3) 一致する場合

 [] は、候補複数識別コード
 [] 内は候補文字
 ↓
 文書認識データ: a b c d e f g N E [B C] o p q r s t u
 文字の一致 : ↑↑ ↑↑
 O O x O ... 一致成立

(4) 一致する場合

 * は、候補多数識別コード
 ↓
 文書認識データ: a b c d e f g N E * o p q r s t u
 文字の一致 : ↑↑↑
 O O O ... 一致成立

【図5】



フロントページの続き

(51)Int.Cl.⁶

識別記号

庁内整理番号

FI

G06F 15/403

技術表示箇所

350C